

07 – reinforcement learning



**ASSIGN
BUSTER**

RL is mapping of X to Y in order to maximise Z. What are X, Y & Z? Actions, situations, reward signal

What does RL explicitly consider? The whole problem of a goal-directed agent interacting with an unknown environment

ON07 - REINFORCEMENT LEARNING SPECIFICALLY FOR YOU FOR ONLY \$13.90/PAGE Order Now

What is a less formal way of describing RL? Learning through trial & error

In RL, what are the two types of feedback the environment provides the agent? State & state evaluation (reward)

In RL, what signal does the agent send to the environment? Action

In RL what does the reward signal communicate? What is required, not how to do it

In RL what does the reward signal tell the agent about the action's correctness? Does not explicitly indicate whether action was correct or incorrect

What are the 3 main differences in supervised learning vs reinforced learning?

- 1.) The system learns from examples by a knowledgeable external factor
- 2.) Environment explicitly indicates what the agent's action should have been
- 3.) Instructive feedback independent of output

What are the 2 main situations where a RL system of learning is appropriate?

- 1.) When it is impossible to get sufficient examples of desired behaviour
- 2.) Learning from experience becomes more appropriate when it becomes difficult for the examples of desired behaviour correct and representative of all situations the agent is likely to experience

What are the 3 main components of an RL algorithm?

- 1.) Reward function
- 2.) Value function
- 3.) Policy

What is the "reward function"? Function that defines a goal by specifying a number for each state-action combination

What is the "value function"? Function that specifies total reward expected when starting from a given state with a given behaviour

What is the "policy"? Mapping from perceived state to action

What is "discounting"? Reward is discounted over longer runs according to some

<https://assignbuster.com/07-reinforcement-learning/>

discount rate with range 0 - 1 What is the discounted return R at time t

(discount function)? (PHOTO) $R_{\nabla t} = r_{\nabla t+1} + \gamma r_{\nabla t+2} + \gamma^2 r_{\nabla t+3} \dots$

$= \sum \gamma^n r_{\nabla t+n}$ What is the action value estimation function (for small k values

where maintaining a prior reward list is feasible)? (PHOTO) $Q_{\nabla t}(a) = (r_{\nabla 1} +$

$r_{\nabla 2} \dots + r_{\nabla k}) / k$ What is the action value estimation function (for large k

values where maintaining a prior reward list is impossible)? (PHOTO) $Q_{\nabla k+1}$

$= Q_{\nabla k} + 1/(k+1) \times (r_{\nabla k+1} - Q_{\nabla k})$

or replace $1/(k+1)$ with a constant for dynamic tasks Describe a greedy

reward policy Highest action value used to select output for given

situation Name 2 other reward policies besides " greedy" ϵ -greedy,

annealed What two things does an environment with the Markov Property

allow us to predict? Its one-step dynamics enable next-state predictions &

expected next reward Compare policy for a Markov state vs. policy as a

function of complete histories They are the same Explain partially-observable

states (the perceptual aliasing problem) If an entity's inputs convey partial

information about the environment, there may be situations which appear

identical to the agent but require different optimal actions SARSA

PSEUDOCODE PHOTOSARSA PSEUDOCODE PHOTO Explain " on-policy" An on-

policy algorithm evaluates the policy actually used Explain " off-policy" An off-

policy algorithm approximates optimal action-value function independently

of the policy being followed ONE-STEP Q-LEARNING FORMULA PHOTO ONE-

STEP Q-LEARNING FORMULA PHOTO What are the 6 main steps in Q-Learning

Pseudocode (PHOTO) 1.) Initialise action value function for states s and

actions a $Q(s, a)$ 2.) For each episode... Initialise s 3.) For each step of trial...

Choose a from s via policy derived from Q 4.) Take action a , observe r, s' 5.)

Update $Q(s, a)$, 6.) $s < -s'$ END foreach END foreach UNTIL terminal Why is

<https://assignbuster.com/07-reinforcement-learning/>

creating a table value for every state-action mapping not feasible in some problems? In complex problems, the data memory required or time to visit all combinations would be too large. How do you generalise in complex problems? Approximate for states not experienced, through supervised learning. What are two ways you can approximate functions during generalisation of complex problems? Gradient descent, MLPs. Briefly describe how you would use an MLP to approximate functions in complex problem generalisation. Use one MLP per action, use the state as the input, MLP returns $Q(s, a)$ per step. Give two applications of RL in game-playing: Backgammon, draughts. Give two applications of RL in engineering: Robotics, lift-allocation. Give two applications on RL in software: Adaptive games, browser agents.